

```

1  -- CODING CONVENTIONS:
2  -- No globals. Line length: 80 chars or less.
3  -- Parse settings from a help string (see top of file).
4  -- This code does no run anything. Rather it is a module to be loaded
5  -- and run by e.g. rlgo.lua)
6
7  -- VARIABLE NAME CONVENTIONS:
8  -- Leading_upper_case : class
9  -- i. : instance var
10 -- l, s : reference to a library function
11 -- prefix _ : some internal function,variable.
12
13 -- TYPE HINT CONVENTIONS (where practical, on function arguments):
14 -- t = table
15 -- prefix s=string
16 -- prefix n=num
17 -- prefix is=boolean
18 -- class names in lower case denote vars of that class
19 -- suffix s denotes table of things
20 local l = require"lib"
21 local the = l.settings[]
22
23 RL.LUA : stings
24 (c)2022 Tim Menzies <tim@ieee.org> BSD (2clause).
25
26 USAGE:
27 lua rlgo.lua [ -bfgkhs [ARG] ]
28
29 OPTIONS:
30 -b --bins discretization control = 8
31 -F --Far in "far", how far to seek = .95
32 -g --go start-up action = pass
33 -h --help show help = false
34 -k --keep keep only these nums = 512
35 -m --Min stop at N*m = .5
36 -p --p distance coefficient = 2
37 -s --seed random number see = 10019
38 -S --Some in "far", how many to search = 512
39 ]]
40 local RL = {About={}, Data={}, Row={}, Col={}, the=the}
41 local About = RL.About -- factory for making columns
42 local Data = RL.Data -- store rows, and their column summaries
43 local Row = RL.Row -- stores one row.
44 local Col = RL.Col -- summarize 1 column. Has 2 roles-- NOMinal,RATIO for syms,nums
45
46 -- FYI: I considered splitting Col into two (one for NOMinals and one for
47 -- RATIOs). But as shown in Col (below), one of those two cases can usually be
48 -- handled as a one-liner. So the benefits of that reorg is not large.
49
50 -- One nuance here is that, to save memory, Rows are created by the FIRST Data
51 -- that sees a record, then shared across every other clone of the data (e.g.
52 -- when clustering, the super data points to the same Row as the sub-data
53 -- cluster of all the other rows closest to that first Row). Since rows
54 -- maintains a pointer to its creator Data object, that first data Data can be
55 -- used to store information about the entire data spaces (e.g. the max and min
56 -- possible values for each columns). This makes certain functions easier like,
57 -- say, distance).
58
59 -- ABOUT
60
61 -- Factory for making columns.
62 function About.new(sNames)
63 return About._cols({names=sNames, all={}, x={}, y={}, klass=nil, sNames} end
64
65 -- How to recognize different column types
66 local is={
67 nom = "[a-z]", -- ratio cols start with uppercase
68 goal = "[!a-z]", -- !klass, {y,-}=maximize,minimize
69 klass = "[*]", -- klass if "*"
70 skip = "$", -- skip if "$"
71 less = "-$]" -- minimize if "-"
72
73 -- Turn a list of column names into Col objects. If the new col is independent
74 -- or dependent or a goal attribute then remember that in i.x or i.y or i.klass.
75 function About._cols(i, sNames)
76 for at,name in pairs(sNames) do
77 local col = l.push(i.all, Col.new(name,at))
78 if not name:find(is.skip) then
79 l.push(name:find(is.goal) and i.y or i.x, col)
80 if name:find(is.klass) then i.klass=col end end end
81 return i end
82
83 -- Update, only the non-skipped cols (i.e. those found in i.x and j.x.
84 function About.add(i,t)
85 local row = t.cells and t or Row.new(i, t)
86 for _,cols in pairs(i.x,i.y) do
87 for _,col in pairs(cols) do
88 Col.add(col, row.cells[col.at]) end end
89 return row end
90
91 -- ROW
92
93 -- Hold one record
94 function Row.new(about,t)
95 return {about=about, cells=t, cooked=l.map(t,l.same)} end
96
97 -- Everything in rows, sorted by distance to i.
98 function Row.around(i,rows)
99 local fun = function(j) return {row=j, d=Row.dist(i,j)} end
100 return l.sort(l.map(rows, fun), l.lt"d") end
101
102 -- Recommend sorting i before j (since i is better).
103
104 function Row.better(i,j)
105 i.evald, j.evald = true,true
106 local s1,s2,d,n,x,y=0,0,0,0,0
107 local ys,e = i._about.y,math.exp(1)
108 for _,col in pairs(ys) do
109 x,y = i.cells[col.at], j.cells[col.at]
110 x,y = Col.norm(col,x), Col.norm(col,y)
111 s1 = s1 - e*(col.w * (x-y)/#ys)
112 s2 = s2 - e*(col.w * (y-x)/#ys)
113 end
114 return s1/#ys < s2/#ys end
115
116 -- Distance
117 function Row.dist(i,j)
118 local d,j,x,y,dist=0,0
119 local cols = cols or i._about.x
120 for _,col in pairs(cols) do
121 x,y = i.cells[col.at], j.cells[col.at]
122 d = d + Col.dist(col,x,y)^the.p
123 n = n + 1 end
124 return (d/n)^(1/the.p) end

```

```

126 -- COL
127
128 -- Summarize one column.
129 function Col.new(txt,at)
130 txt = txt or ""
131 return {n = 0, -- how many items seen?
132 at = at or 0, -- position of column
133 txt = txt, -- column header
134 isNom=txt:find(is.nom),
135 w = txt:find(is.less) and -1 or 1,
136 ok = true, -- false if some update needed
137 _has = {}} end
138
139 -- Create columns with particular roles.
140 function Col.ratio(...) local i=Col.new(...); i.isNom=false; return i end
141 function Col.nom(...) local i=Col.new(...); i.isNom=true; return i end
142
143 -- Update. Optically, repeat n times.
144 function Col.add(i,x, n)
145 if x ~= "" then
146 n = n or 1
147 i.n = i.n + n
148 if i.isNom then i._has[x] = n + (i._has[x] or 0) else
149 for = 1,n do
150 local pos
151 if #i._has < the.keep then pos= 1 + (#i._has)
152 elseif l.rand() < the.keep/i.n then pos=l.rand(#i._has) end
153 if pos then
154 i.ok=false -- kept items are no longer sorted
155 i._has[pos]=x end end end end
156
157 -- Distance. If missing values, assume max distance.
158 function Col.dist(i,x,y)
159 if x=="?" and y=="?" then return 1 end
160 if i.isNom then return x==y and 0 or 1 else
161 if x=="?" then y = Col.norm(i,y); x=y<.5 and 1 or 0
162 elseif y=="?" then x = Col.norm(i,x); y=x<.5 and 1 or 0
163 else x,y = Col.norm(i,x), Col.norm(i,y) end
164 return math.abs(x-y) end end
165
166 -- Diversity: divergence from central tendency (sd,entropy for NOM,RATIO).
167 function Col.div(i)
168 local t = Col.has(i)
169 if not i.isNom then return (l.per(t,.9) - l.per(t,.1))/2.58 else
170 local e=0
171 for _,v in pairs(t) do if v>0 then e=-v/i.n*math.log(v/i.n,2) end end
172 return e end end
173
174 -- Sorted contents
175 function Col.has(i)
176 if i.isNom then return i._has else
177 if not i.ok then table.sort(i._has) end
178 i.ok=true
179 return i._has end end
180
181 -- Central tendency (mode,median for NOMs,RATIOs)
182 function Col.mid(i)
183 if not i.isNom then return l.per(Col.has(i),.5) else
184 local mode,mst=nil,-1
185 for k,v in pairs(i._has) do if v>mst then mode,mst=k,v end end
186 return mode end end
187
188 -- Return num, scaled to 0..1 for lo..hi
189 function Col.norm(i,x)
190 if i.isNom then return x else
191 local has=Col.has(i) -- "a" contains all our numbers, sorted.
192 local lo,hi = has[1],has[#has]
193 return hi - lo < 1E-9 and 0 or (x-lo)/(hi-lo) end end
194
195 -- Map x to a small range of values. For NOMs, x maps to itself.
196 function Col.discretize(i,x)
197 if i.isNom then return x else
198 local has = has(i)
199 local lo,hi = has[1], has[#has]
200 local b = (hi - lo)/the.bins
201 return hi=lo and 1 or math.floor(x/b+.5)*b end end

```

```

204 -- DATA
205
206 -- Holds n records
207 function Data.new(t) return {rows={}, about=About.new(t) } end
208
209 -- Update
210 function Data.add(i,t) l.push(i.rows, About.add(i.about,t)) end
211
212 -- Sort rows, then pretend you didn't
213 function Data.cheat(i)
214 for j,row in pairs(l.sort(i.rows, Row.better)) do
215 row.rank = l.rand(100*j/#i.rows)
216 row.evald=false end
217 i.rows = l.shuffle(i.rows) end
218
219 -- Replicate structure
220 function Data.clone(i, t)
221 local out = Data.new(i.about.names)
222 for _,row in pairs(t or {}) do Data.add(out,row) end
223 return out end
224
225 -- Discretize all row values (writing those vals to "cooked").
226 function Data.discretize(i)
227 for _,row in pairs(i.rows) do
228 for _,col in pairs(i.about.x) do
229 local x = row.cells[col.at]
230 if x=="?" then
231 row.cooked[col.at] = Col.discretize(col,x) end end end end
232
233 -- Diversity
234 function Data.div(i) return l.map(i.about.y, Col.div) end
235
236 -- Recursively bi-cluster one Data into sub-Datas.
237 function Data.cluster(i, rowAbove,stop)
238 stop = stop or (#i.rows)*the.Min
239 if #i.rows >= 2*stop then
240 local A,B,As,Bs,c = Data.half(i.rows,rowAbove)
241 i.halves = {c=c, A=A, B=B,
242 kids = { Data.cluster(Data.clone(i,As), A, stop),
243 Data.cluster(Data.clone(i,Bs), B, stop) }}end
244 return i end
245
246 -- Split data according to distance to two distant points A,B
247 -- To dodge outliers, don't search all the way to edge (see the.Far).
248 -- To speed things up:
249 -- - try to reuse a distant point from above (see rowAbove).
250 -- - only look at some of the rows (see the.Some).
251 -- - find distant points in linear time via
252 -- A=far(any(i)) and B=far(A).
253 function Data.half(i, rows, rowAbove, c)
254 local some= l.many(rows, the.Some)
255 local function far(row)
256 return l.per(Row.around(row,some), the.Far).row end
257 local As,Bs = {},{}
258 local A= rowAbove or far(l.any(some))
259 local B= far(A)
260 local c= Row.dist(A,B)
261 local function project(row)
262 local a,b = Row.dist(row,A), Row.dist(row,B)
263 return {row=row, x=(a^2 + c^2 - b^2)/(2*c)} end
264 for n,row in pairs(l.sort(l.map(rows, project), l.lt"x")) do
265 l.push(n < #rows/2 and As or Bs, rowx.row) end
266 return A,B,As,Bs,c end
267
268 -- Load from file
269 function Data.load(sFilename, data)
270 l.csv(sFilename, function(row)
271 if data then Data.add(data,row) else data=Data.new(row) end end)
272 return data end
273
274 -- Central tendency
275 function Data.mid(i)
276 local t={}
277 for _,col in pairs(i.about.y) do t.n=#i.rows; t[col.txt] = Col.mid(col) end
278 return t end
279
280 function Data.best(i, rowAbove,stop)
281 stop = stop or 10 --(#i.rows)*the.Min
282 if #i.rows < stop
283 then return i.rows
284 else local A,B,As,Bs,c = Data.half(i, i.rows, rowAbove)
285 if Row.better(A,B)
286 then return Data.best(Data.clone(i,As),B,stop) end end end
287
288 -- Heuristically sort rows by trends in the y-values
289 -- (specified, evaluated two remote points, sort worse
290 -- half by distance to best point, recurse on best half)
291 function Data.trends(i,out, rowAbove,stop)
292 stop = stop or (#i.rows)*the.Min
293 out = out or {}
294 if #i.rows < stop
295 then for _,row in pairs(i.rows) do l.push(out,row) end
296 else local A,B,As,Bs,c = Data.half(i, i.rows, rowAbove)
297 if Row.better(A,B)
298 then for j=#Bs,1,-1 do l.push(out,Bs[j]) end
299 Data.trends(Data.clone(i,l.rev(As)), out,A, stop)
300 else for _,row in pairs(As) do l.push(out,row) end
301 Data.trends(Data.clone(i,Bs), out,B, stop) end end
302 return out end
303
304 -----
305 return RL

```